END
DATE
FILMED

1  84

DTIC

MICROCOPY RESOLUTION TEST CHART

NATIONAL BUREAU OF STANDARDS 1963 A

# AN ELEMENTARY EXPOSITION OF THE HARSANYI—SELTEN TRACING PROCEDURE

## LEVEL Ⅱ ⑫

Final Report

November 1980

DTIC
S ELECTE
DEC 0 9 1980
E

Prepared for:

DDC FILE COPY

SRI International

80 12 08 030

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER | 2. GOVT ACCESSION NO. AD-A092 687 | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle) AN ELEMENTARY EXPOSITION OF THE HARSANYI-SELTEN TRACING PROCEDURE . | | 5. TYPE OF REPORT & PERIOD COVERED Technical Report |
| | | 6. PERFORMING ORG. REPORT NUMBER 7813 |
| 7. AUTHOR(s) Horace W. Brock, Project Leader | | 8. CONTRACT OR GRANT NUMBER(s) N00014-78-C-0731 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS SRI International 333 Ravenswood Avenue Menlo Park, California 94025 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS |
| 11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research 800 N. Quincy Street Arlington, Virginia 22217 | | 12. REPORT DATE November 1980 — 13. NO. OF PAGES 58 |
| | | 15. SECURITY CLASS. (of this report) UNCLASSIFIED |
| 14. MONITORING AGENCY NAME & ADDRESS (if diff. from Controlling Office) | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE N/A |

16. DISTRIBUTION STATEMENT (of this report)

Unlimited

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from report)

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)
Game Theory
Equilibrium Points
Noncooperation Solution Theory
Tracing Procedure

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

J. Harsanyi and R. Selten have recently developed a new theory that attempts to
determine a unique solution to any game. In Chapter 1 the historical and
intellectual background of the new Harsanyi-Selten "tracing procedure" is given.
In Chapter 2, the author presents an elementary exposition of the tracing
procedure itself. Finally in CHapter 3, the author presents a worked sample
application of the tracing procedure in the form of an analysis of a
one parameter family of two-person bargaining games. Throughout, there is an
emphasis on the intuitive aspects of the new theory.

**DD** FORM 1473
1 JAN 73
EDITION OF 1 NOV 65 IS OBSOLETE

CONTENTS

## 1. INTRODUCTION AND OVERVIEW

During the past eight years, John Harsanyi and Reinhard Selten
have collaborated on the development of a new approach to game theory.
This approach centers upon an operation they call the "tracing
procedure". Their new theory is at once very complex, very general
and very ambitious: it attempts nothing less than to prescribe a
unique solution to any and every competitive situation that can pro-
perly be called a "game".

In this research paper, we attempt to describe the tracing procedure
and to motivate it by applying it to a simple worked example. However,
before proceeding to the intricacies of the Harsanyi-Selten theory, we
shall attempt here in Chapter 1 to put the new research effort in proper
perspective. We ask and partially answer such queries  as: Why have
Harsanyi-Selten been prompted to develop the theory? What predecessor
theories were there, and what was wrong with them? Is the new theory
likely to be useful in practical applications? In Section 1.1 we
discuss two foundational problems in game theory that prompted the new
theory. In 1.2 we first (Section 1.2.A) discuss an earlier theory that
Harsanyi developed and then abandoned, the so-called bargaining equilibrium
model for non-cooperative games. Then in 1.2.B we contrast the new tracing
theory with the older theory. In Section 1.3 we argue that the new theory
is not likely to be useful in practical applications, for reasons having
to do not only with the theory but with the purpose of game theory itself.

Subsequently, in Chapter 2, we offer a somewhat technical description
of the tracing procedure itself. This second chapter also contains
a philosophically oriented appraisal of the procedure that emphasizes
its decision theoretical and game theoretical interpretations. Then
in Chapter 3 we furnish a worked application analyzed in recent
unpublished work by Harsanyi (and also by Guth and Selten). A
two-person bargaining game with incomplete information is constructed.
It is shown that the tracing procedure yields a solution that is
almost identical to the classical two-person Nash bargaining solution
of a game of complete information that is almost identical to the example
being analyzed. This result is of interest not only because of the
historical interest and importance of the Nash bargaining solution,
but also because of the well-known intuitive appeal of the solution.

## 1.1. Two Foundational Problems in the Theory of Games

While there are several fundamental "difficulties" with game theory
as a mathematical discipline, two problems stand out as fundamental.
First, there is the somewhat arbitrary and obfuscatory division of
competitive situations into two different classes: "cooperative"
and "non-cooperative". Classical game theory follows Nash in defining
cooperative games as those permitting both free communication between
the players and also enforceable agreements between them. All other
games are non-cooperative, although classical game theory usually
assumes that a non-cooperative game is characterized by an absence

both of binding agreements and communication. The second fundamental
problem is that of defining a unique solution to games of various types.
Concerning solution theory, the ideal has been to do what von Neumann
originally did in his original 1928 paper on zero-sum games: namely
to characterize a unique solution for all games of the class under
consideration. Regrettably, with the exception of the class of zero-sum
games and certain two-person cooperative games (e.g. the Nash bargaining
games), there are very few classes of games for which determinate
solution theories have been formulated. And to the extent that determinate
theories have been set forth (e.g. the Shapley value for n-person
transferable utility games) there is an even deeper problem: Why should
different classes of games require fundamentally different solution
theories, each with a different "logic"? Wouldn't it be preferable
to have a single overarching solution theory that would be based on a
single set of principles applying to all games whatsoever?

## 1.2. The New Harsanyi-Selten Theory

The new Harsanyi-Selten theory attempts to solve both of the foundational
problems cited above in one fell swoop. It takes the concept of
a non-cooperative game as basic, and essentially requires that
any given cooperative game to be analyzed be reformulated as
a non-cooperative game. The theory then restricts its attention to
the set of equilibrium points of the non-cooperative game, and

attempts to select one of these equilibrium points as the solution
to the game. Exactly how this selection process operates is what
we shall be discussing in subsequent chapters. In the present chapter,
we shall merely sketch the spirit of the process. Before doing this,
however, it will be helpful to review the only other approach to
defining a unique solution to all games -- John Harsanyi's "Bargaining
Equilibrium" theory set forth in 1977 in Harsanyi ( 5).[*]

### 1.2.A. The Bargaining Equilibrium Solution Theory

Harsanyi's earlier work will be of interest to us not only because
his theory was the first theory to attempt to define a unique solution
to all games based on a single set of principles, but also because
it made use of a concept ("risk dominance") that plays a modified
role in the newer Harsanyi-Selten theory.

In 1950, John F. Nash, Jr. ( 8) had axiomatized a unique solution
to the two-person bargaining game with complete information. While
Nash's model was mathematically very elegant and abstract, it failed
to provide any psychologically convincing reasons as to why rational
bargaining in the real world would result in the "solution" prescribed
by Nash's theory. Harsanyi ( 2) remedied this deficiency in an

[*] While H͏ ͏yi did not publish his bargaining equilibrium approach
in book fo͏ until 1977, his work was largely completed by 1970.
And it was around 1970 that Harsanyi joined Selten in the newer
Harsanyi-Selten theory that is being discussed in the present paper.

important 1956 paper. Building upon and extending an earlier model
of F. Zeuthen (11), he views bargaining as a dynamic process in
which the antagonists make "demands" and "concessions". Indeed, suppose
that we are at stage k in a two-person bargaining game. Let $U_1$
and $U_2$ be the bargainers' respective utility functions. Denote their
last offers at stage k as $A_1^k = A_1$ and $A_2^k = A_2$. The critical question
is: Which player will have to make the next concession at stage k+1?
The Harsanyi-Zeuthen argument is that the next concession must come
from the player who is less willing to face the risk of a conflict.
(In the model, the "conflict situation", denoted by C, represents what
happens if negotiations collapse and both players go away with nothing.)
But how can we measure a given player's willingness to risk a conflict
rather than to accept an opponent's terms? In explaining how the
Harsanyi-Zeuthen theory answers this critical question, I shall
follow Harsanyi ( 5 , pp. 150-155).

Because we want to measure each player's willingness to stick to his
own terms rather than to accept his antagonist's terms, we can view
player i's decision as a choice between two alternatives, a' and b'.
Here a' denotes full insistence on i's own last offer $A_i$, and b'
denotes full acceptance of his opponent's last offer $A_j$. Assuming
i is "rational" (i.e. that he is a Bayesian expected utility maximizer)
then he must start by assigning subjective probabilities to the two

possible choices that his opponent can make in reply. Let $p_{ji}$ be the subjective probability that i assigns to the hypothesis that j will chose alternative a'; and let $q_{ji} = 1 - p_{ji}$ denote the subjective probability that i assigns to the hypothesis that j will choose b'.

If i himself chooses b', then according to the Harsanyi-Zeuthen model, he will obtain payoff $U_i(A_j)$ with certainty, regardless of what j does. On the other hand, if i chooses alternative a', then he may obtain the higher payoff $U_i(A_i) > U_i(A_j)$; but he may also obtain the conflict payoff $U_i(C) < U_i(A_j)$. The former outcome will occur with probability $q_{ji}$ whereas the latter will occur with probability $p_{ji}$.

Thus if player i wants to maximize his expected utility, then he will stand on his own last offer $A_i$, that is select a', only if

$$(1 - p_{ji}) \cdot U_i(A_i) + p_{ji} \cdot U_i(C) \gtrless U_i(A_j) \tag{1.1}$$

Equivalently, i will select a' only if

$$p_{ji} \lessgtr r_i = \frac{U_i(A_i) - U_i(A_j)}{U_i(A_i) - U_i(C)} \tag{1.2}$$

The quantity $r_i$ (i = 1,2) defined in (1.2) is called player i's risk limit since it represents the highest risk (the highest subjective probability of a conflict) that player i would be willing to face in order to achieve an agreement on his own terms $A_i$ (rather than

on his opponent's terms $A_j$). This is so because, if player i sticks
to his own last offer $A_i$, then he must expect a conflict to occur
with probability $p_{ji}$ (this is the probability that his opponent will
stick to his own last offer $A_j$, in which event neither makes a concession
and a conflict results). But, from (1.2) we know that the highest value
of the probability $p_{ji}$ that i can face, without _switching over_ to
acceptance of the opponent's last offer, is $p_{ji} = r_i$.

We have set forth the basic logic of the Harsanyi-Zeuthen scheme.
Several important points should be added to the foregoing sketch of
their model. First, Harsanyi shows that each player i will always
know the subjective probabilities $p_{ji}$ and $q_{ji}$ he needs in performing
his "utility maximization" according to the above logic. We do not have
space to show why this is so, but it is important since it offers one
of the few cases where "subjective" probabilities are in fact quasi-
"objective". Second, the Harsanyi-Zeuthen model implies that the
bargaining process will stop -- and agreement will be reached -- at
that point where the arithmetic product of the players' utilities is
reached. But this is the well-known prescription of Nash's abstract
model. Hence the two models imply the same bargaining "solution",
namely the set of terms whereby the product-of-the-utilities is
maximized. We shall not work out the details of this equivalence
here, but refer the reader to the book by Harsanyi.

Third, note that we can apply our analysis to j as well as to i, and we can compute the quantity $r_j$ (as well as $r_i$) when evaluating any set of bargaining terms. Thus, we can define a <u>risk dominance</u> concept along the following lines. If $r_i < r_j$, then we know that player i is less willing than player j to risk a conflict, or equivalently, that he has less incentive to do so. But both players will know this. Therefore player i will be under strong <u>psychological</u> <u>pressure</u> to make the next concession, whereas player j will be unwilling to do so. This is the essence of the concept of risk dominance that in one form or another plays an important role both in Harsanyi's first general theory (described just below) of 1977, and in the later Harsanyi-Selten theory. We shall now briefly discuss the essence of Harsanyi's "bargaining equilibrium approach" to defining a unique solution to all games.

In discussing the problem of multiple equilibrium points and the need for a determinate solution theory in game theory, Luce and Raiffa ( 7, pp. 109-110) present the following non-cooperative game:

|        | $B_1$     | $B_2$    |
|--------|-----------|----------|
| $A_1$  | (4, -30)  | (10, 6)  |
| $A_2$  | (8,8)     | (5,4)    |

- 9 -

Here player 1 has a choice between strategies $A_1$ and $A_2$. Player 2 has a choice between $B_1$ and $B_2$. The two pairs of strategies $(A_2,B_1)$ and $(A_1,B_2)$ are both equilibrium points of the game: that is, (i) $A_2$ is a best reply strategy to $B_1$ just as $B_1$ is a best reply to $A_2$, and (ii) $A_1$ is optimal against $B_2$ just as $B_2$ is optimal against $A_1$. We see here an instance of perhaps the most fundamental problem in game theory, that of the multiplicity of equilibrium points in non-cooperative games. In the present case, is there some way of saying that $(A_1,B_2)$ is "better than" or "worse than" $(A_2,B_1)$?

Luce and Raiffa argue as follows. If player 2 has any reason to fear that 1 will take $A_1$, then he dare not take $B_1$ for fear of getting -30. But player 1, knowing this, has every reason to take $A_1$ which gives him his maximum payoff. But now the argument is cyclical. For 2, having some rationalization for 1's adoption of $A_1$ has all the more reason to avoid $B_1$. Thus the equilibrium point $(A_1,B_2)$ psychologically dominates $(A_2,B_1)$.

Luce and Raiffa's suggestion here is merely a suggestion, and a rather imprecise one at that. They have not really made clear either when psychological dominance will occur (in some games there may be no psychologically dominating equilibrium points) or what the exact nature of this dominance really is.

What Harsanyi did in his first general solution theory (his "bargaining equilibrium approach") is to show that (i) the Harsanyi-Zeuthen concept of risk dominance originating in two-person cooperative bargaining games can be used to give a precise quantitative meaning to Luce and Raiffa's "psychological dominance" concept; and (ii) the risk dominance concept can provide the basis for a determinate solution theory that identifies a unique equilibrium point as the solution to any given non-cooperative game. To see how Harsanyi's first theory works, let us consider the following simple example set forth by Harsanyi ( 5, pp. 275-276). Consider the following non-cooperative game:

|       | $B_1$     | $B_2$      |
|-------|-----------|------------|
| $A_1$ | (2,1)     | (-10,-1)   |
| $A_2$ | (0,0)     | (1, 2)     |

In this game there are only two equilibrium points, namely $s^1 = (A_1, B_1)$ and $s^2 = (A_2, B_2)^*$. Obviously player 1 will prefer $s^1$ whereas player 2 will prefer $s^2$. If both players use their strategies associated with their own preferred equilibrium points, then player 1 would suffer a

---

* There is a third equilibrium point, namely (1/2 $A_1$ + 1/2$A_2$, 11/13$B_1$ + 2/13$B_2$). But for rather technical reasons, Harsanyi rules this equilibrium point out of consideration in his theory. It is "ineligible".

much greater loss (in relation to the payoff difference for him
between the two equilibrium points) than player 2: for $U_1(A_1,B_2)$
= -10, while $U_2(A_1,B_2)$ = -1. Hence both players will know that
player 1 will be more fearful to use strategy $A_1$ in defiance of
player 2's using $B_2$ than player 2 will be fearful to use $B_2$ in
defiance of player 1's using $A_1$. Thus we see that, in Luce and
Raiffa's terminology, both players will settle down at equilibrium
point $s^2 = (A_2,B_2)$ preferred by player 2 since this point psychol-
ogically dominates point $s^1$.

Harsanyi restates this point in terms of risk dominance as follows.
He views the choice between the two eligible equilibrium points as
a kind of tacit bargaining game between the two players. Assuming
each player attempts to maximize his expected utility, then, for
reasons we have seen, player 1 will stick to strategy $A_1$ as long as
the subjective probability that he attaches to the hypothesis that
player 2 will use strategy $B_2$ is no greater than

$$r_1 = \frac{U_1(A_1,B_1)- U_1(A_2,B_2)}{U_1(A_1,B_1)- U_1(A_1,B_2)} = \frac{1}{12} = .08 \qquad (1.3)$$

Likewise player 2 will stick to strategy $B_2$ as long as the subjective
probability that he attaches to the hypothesis that player 1
will stick to $A_1$ is no greater than his risk limit $r_2$. $r_2$ is
easily calculated to be .33. By the Harsanyi-Zeuthen concession

principle, since $r_1 = .08$ is smaller than $r_2 = .33$, player 1 must
yield and accept as the solution to the game the equilibrium point
$s^2 = (A_2, B_2)$ preferred by his opponent. In brief, we say that $s^2$
<u>risk</u> <u>dominates</u> $s^1$.

Harsanyi's general theory for defining a solution to any n-person
non-cooperative game is an elaborate generalization of the ideas
we have seen worked out in the above example. As we have argued
in a recent review of Harsanyi's 1977 book (5, pp. 358-361), this
bargaining equilibrium solution theory is only a partial success.
There are two basic reasons for this judgment, and these views are
shared by Harsanyi. First, the theoiy is too complex. If the theory
boiled down simply to establishing a linear ordering of the different
eligible equilibrium points of a given game by using the generalized
risk dominance relations Harsanyi develops, then all would be well.
But this is not the case. In order to get a determinate theory, Harsanyi
finds himself to classify the various equilibrium points of a game
as "weak", "strong", "centroid", "weakly/strongly efficient", "inefficient",
"admissible", and so on. It is all too complex. Second, there are
cases where the players in Harsanyi's theory will end up using strategies
that do not constitute an equilibrium point. They are obliged
rather to use maximin strategies. This is unsatisfactory since, as
Harsanyi usually insists, any "rational" solution to any game must
at a minimum be an equilibrium point of the game.

### 1.2.B. The New Harsanyi-Selten Theory

Prompted in part by the criticisms we have just advanced of Harsanyi's bargaining equilibrium solution theory, Harsanyi and Selten proceeded to develop a different determinate solution theory. Since an operation called the "tracing procedure" lies at the heart of this theory, the entire theory is often spoken of as the "tracing procedure theory".

Since we are about to examine the tracing procedure in considerable detail, the only thing we shall attempt here is to suggest how this new theory contrasts with the older bargaining equilibrium model of Harsanyi. As we see it, there are two distinctive differences between the two approaches. First, the bargaining equilibrium approach has a distinctive cooperative game spirit about it. To begin with, the games analyzed by Nash and by Harsanyi-Zeuthen in their papers on the bargaining problem were cooperative games. Moreover, the Harsanyi-Zeuthen model of bargaining as a process of demands and concessions resulting in a binding agreement is far more suitable to cooperative games (in which binding agreements and communication are possible) than to non-cooperative games. In contrast to this, the tracing procedure theory is distinctively non-cooperative in spirit. According to this theory, any game whatsoever that is to be "solved" is reduced to a non-cooperative game. Thereafter, the solution theory operates solely on the equilibrium points of the game. And in selecting one equilibrium point as the "solution" to the game, the tracing procedure

does not rely upon a principle of "tacit bargaining" or anything else that smacks of classical cooperative games.

The second fundamental difference between the old and the new solution theory lies in the psychological model that is developed to yield a determinate solution. The bargaining equilibrium model makes use of the concept of psychological dominance, or more specifically risk dominance. As we have seen, symmetrically rational players will, through a process of tacit bargaining, end up playing the strategies of that equilibrium point that risk dominates all other equilibrium points. In the tracing procedure theory, the model that brings symmetrically rational players into agreement as to which equilibrium point strategies to play is both cognitive and psychological. At this point in this essay, it is difficult to explain exactly what we mean by this assertion. Briefly, during the operation of the tracing procedure, the players start off by assigning prior probabilities to each other's behavior. These prior probabilities are based on what Harsanyi and Selten call "naive Bayesian rationality". Specifically, the theory makes the prior probability that player i will use strategy $s_i$ as opposed to $s_i'$ proportional to his net payoffs associated with the two equilibrium points s and s'. This form of rationality is "naive" because it turns out that the strategies recommended to each player via this approach do not constitute an equilibrium point of the

underlying game being analyzed. But a "solution" must (according to
game theory) be an equilibrium point. Thus it is necessary for the
players to "modify" the naive behavior in such a way that they
end up agreeing on a genuine equilibrium point as the solution.
But how do they modify their initial, naive behavior? They do so
by viewing the information contained in their initial strategy
inclinations as "prior" information that is "updated" according to
the rather complex mechanics of the tracing procedure itself (cf.
Chapter 2). When this updating is completed, a complete convergence
of strategic intentions and expectations has been achieved, and all
the players tacitly agree which equilibrium point is the solution.
During this process of updating, we shall observe that certain
players "switch" from their initial (naive) strategies to the
strategies corresponding to other equilibrium points. The reason
why switching takes place, and the order in which the various players
do (or do not) switch, is explained in terms of a new concept of
risk dominance that Harsanyi-Selten introduce. We shall observe this
in detail in Chapter 3.

Thus, the tracing procedure solution theory contains a rich cognitive
model concerning the manner in which the players in a game form
preliminary expectations about each other's behavior, and subsequently
modify these expectations via informational updating. But the
updating process relies in part upon a rather psychological model of

risk dominance instantiated by the <u>strategy-switching behavior of</u> the players as they grope their way in the dark to a "solution" that is tacitly agreed upon by all. Thus it is that the new theory differs from the bargaining equilibrium theory in containing a strong cognitive component that is lacking in the former. All this will become more clear in the sequel.

Hopefully, the foregoing remarks serve to place in a rather broad perspective the relationship bet' en the new Harsanyi-Selten theory and the only other previous attempt to define a determinate solution theory applicable to all games.

## 1.3 The Applicability of the New Theory to Real-World Problems

One question that immediately arises when discussing such a rich, general and complex theory as the Harsanyi-Selten theory is whether it is likely to be of "practical" interest. Our present and rather tentative belief is that the Harsanyi-Selten theory is not likely to prove to be of practical importance. The reason for this lies partly in the nature of the new theory itself, and partly in the nature and purpose of game theory as a whole. Let us begin with the second of these two reasons.

It is frequently asserted that game theory assumes that people are perfectly rational. But as Oskar Morgenstern pointed out to the present writer on several occasions, this is a misconception.

The fundamental purpose of game theory -- as well as its greatest

contribution to date -- lies in its success in making clear what

it means for people to behave "rationally" in game situations. The

theory is prescriptive and philosophical in nature. But herein lies

one of the limitations of game theory as regards practical applications.

People are not in fact perfectly rational. Hence game theory is not

particularly useful for predicting what people will do in a given

competitive situation. Let us explore this difficulty at a somewhat

deeper level by contrasting the potential usefulness of decision theory

on the one hand, and game theory on the other hand.

Suppose there are n players who are involved in a competitive situation,

loosely, a "game". Let us take the vantage point of player i in

deciding whether or not to use game theory, decision theory, or

neither in determining a best course of action. ( Player i might

denote the United States in a world of n military states that interact).

If  i  uses decision theory, he will be able to determine a strategy

that is rational for him  regardless of whether he believes any

or all of his antagonists to be rational. More pertinently, if he

has reason to suppose that any or (more typically) all of his

antagonists are irrational, he can bring this information to bear

on his own strategy selection problem in the following manner. He

can use the information when assigning his prior probabilities to

the actions of his antagonists. Then, given these priors, and given

any other (structural) information that he possesses, he can use
decision theory to determine a course of action that is optimal
for himself. Thus, decision theory is highly flexible: it can be
used by party i in a wide variety of cases, most importantly in
those cases where he has reason to believe that his antagonists are
not perfectly rational in a game theoretical sense. Note that in
stating this we are not saying how player i can defensibly choose his
prior probabilities concerning his antagonists' action. This is a
notoriously difficult problem, and is one of the problems that
gave rise to the theory of games in the first place.

If player i wishes to use game theory rather than decision theory
he will immediately run into two problems. First, game theory instructs
him <u>not</u> to use game theory unless he is sure that <u>all</u> (n-1) antagonists
are perfectly rational in a game theoretic context. The reason for
this is that game theory is a prescriptive theory of <u>symmetrically</u>
<u>rational</u> <u>behavior</u> that assumes that all players have access to and
follow the rules of "game theory". And here we run into the
second problem. There are many different versions of game theory ;
in particular, there are a host of rival definitions of what constitutes
a "solution". In short, even if all n players wish to act perfectly
rationally and expect each other to do so, they would have no way
of telling which "brand" of game theory each other subscribes to.
The upshot of all this is that game theory is as inflexible and useless

for practical applications as Bayesian decision theory is flexible
and useful.

The second reason why the new Harsanyi-Selten theory is not likely
to prove useful in practice is related to the first reason. Even
if all n players in our fictitious game wished to act perfectly
rationally, it is most unlikely that they would identify perfect
symmetric rationality with Harsanyi-Selten behavior. Moreover, even
if they did all agree with the Harsanyi-Selten characterization of
rationality in game situations, we feel it most unlikely that they
could operationalize the theory in the manner required to actually
compute a solution. Both time constraints as well as the probability
of interpretative error underly this assertion. We suspect that the
reader will come to agree with our position here after reading Chapters
2 and 3 in the sequel where the intricacies of the new theory are
explored in some detail.

In advancing the foregoing criticism of the Harsanyi-Selten theory,
we do not in any way intend to denigrate its ultimate significance.
Indeed, this theory lies in the best tradition of game theory. For it
is a notable philosophical contribution attempting to clarify what
"perfect rationality" means in a very general class of game situations.
Moreover, Harsanyi and Selten have tackled head-on the most difficult
problem of all, namely characterizing a unique solution to all games.

## 1.4. A Note on the Content of Chapters 2 and 3

In the following chapters, we first discuss the theory of the tracing
procedure. Thereafter, we draw on an example due to Harsanyi and
show how the tracing procedure works to define a solution to a
certain family of two-person bargaining games of incomplete information.
Our intent has been to produce a paper that is intelligible by as
large a cross section of readers as possible. Harsanyi has described
his new theory in several articles, e.g. (3 and 4). However, anyone
familiar with these expositions is aware of two problems: first, the
new theory appears as extremely complicated. Second, it is difficult
to penetrate the formal theory because no practical applications
are given. The present treatment differs from any other we know of
because it attempts to overcome both these problems in an effort
to broaden the audience interested in and familiar with the theory.

Specifically, our presentation of the formal theory in Chapter 2
(and partly in Chapter 3) takes certain liberties aimed at
streamlining the theory. We present only those portions of the
formal theory that are needed both to communicate the basic ideas
underlying the tracing procedure and to permit comprehension of
the worked example of Chapter 3. This approach reduces the complexity
of the presentation, permits it to be briefer and more accessible,
and ensures that all the formal ideas presented are relevant to
the practical example. The price paid for all this is an omission

of a discussion of (i) the uniform perturbation of games in agent
normal form (even though we do discuss the agent normal form), (ii) cells
and formations -- concepts that while interesting and important are
quite complex and play no role in the example of Chapter 3, and (iii)
the discretization of the family of games $G(w)$ analyzed in Chapter 3.

With these caveats in mind, let us now turn to the new Harsanyi-Selten
theory and learn what the tracing procedure is, and how it works.

## 2. THE TRACING PROCEDURE

In Chapter 1 we provided both a brief description of and the motivation for the new HS "tracing procedure". In this Chapter we shall give a somewhat technical description of the tracing procedure itself. In terms of demands upon the reader, this Chapter falls midway between Chapter 1 (which is accessible to anyone with an elementary interest in game theory) and Chapter 3 (which is necessarily technical at points).

In 2.1 below, we give a somewhat precise overview of the tracing procedure, emphasizing the question of the prior probability distributions that are so important in the theory. In 2.2 we turn to the tracing procedure itself (viewed as a process that operates on a given n-tuple of prior distributions) and discuss one particular variant of it in detail: the linear tracing procedure. Finally, in 2.3 we give a rather philosophical appraisal of this procedure, emphasizing its decision theoretical and game theoretical interpretations and inspirations.

### 2.1 The Tracing Procedure and the Priors

As mentioned in Chapter 1, the tracing procedure is essentially a mathematical model of a process of convergent expectations entertained by the n players in any game. More specifically, the tracing procedure (henceforth, TP) models a solution process by which rational people will come to adopt, and come to expect each

other to adopt, one specific equilibrium point $q* = (q_1*,\ldots,q_n*)$ as the solution for a given noncooperative game G.

At the beginning of this process, the players are ignorant as to which strategies will be used by their mutual antagonists. Thus, according to generalized Bayesian logic, each player j will express his expectations about the strategy choice of any other player $i \neq j$ in the form of a subjective probability distribution $p_i = (p_{i1},\ldots,p_{iK_i})$ over i's pure strategies, where $p_{ik}$ $(k = 1,\ldots,K_i)$ is the probability that player j assigns to the hypothesis that player i will actually use his kth pure strategy $\phi_i^k$. These subjective distributions $p_i$ are to be called the prior probability distributions or simply the priors of the TP.

One very important aspect of the HS theory is that it prescribes the prior $p_i$ of each and every player in any game G. Additionally, the HS theory goes further than this, and it assumes that all (n-1) antagonists of player i will in fact use this prior $p_i$ prescribed by the theory. The question of "where these priors $p_i$ come from" will best be answered in Chapter 3 where we actually compute a family of priors in the context of a concrete family of games. Briefly, in most cases, the HS theory priors have the property that a given prior $p_i$ will assign higher probability to those pure strategies that would be best replies to the other players' expected strategies in wider ranges of possible strategic situations.

Formally speaking, every prior $p_i$ is a probability distribution over player i's pure strategy set. Thus it has the nature of a "mixed strategy" in game theoretic terms. However its game theoretic interpretation is very different from that of a standard mixed strategy for two specific reasons: (i) the probabilities $q_{ik}$ associated with a true mixed strategy in game theory are <u>objective</u> probabilities chosen by player i himself on the basis of a solution theory, whereas the probabilities $p_{ik}$ associated with an HS prior are <u>subjective</u> probabilities expressing the other players' expectations about player i's likely behavior. (ii) In a true mixed strategy, the probabilities denote the intentions of the players to deliberately randomize their pure strategies. Obviously, this is not the case in the HS theory as described above.

$p = (p_1, \ldots, p_n)$ will henceforth denote the n-tuple of all priors in the game, whereas the (n-1)-tuple obtained from p via deletion of the ith component $p_i$ will be expressed as $p_{-i} = (p_1, \ldots, p_{i-1}, p_{i+1}, \ldots, p_n)$. p will be called a <u>prior vector</u> and $p_{-i}$ will be called an <u>i-incomplete prior vector</u>.

At this point we turn to a problem HS call the "prediction problem". Let us take as given the n players' mutual expectations about one another's behavior as given by the prior vector p. How can we now predict the actual strategy combination $q = (q_1, \ldots, q_n)$ these players will use in practice? The entire purpose of the TP is to answer this problem.

At this point we come to a discussion of the dynamics of the TP itself. We have already sketched how it is possible for the players to formulate their mutual expectations about each others' likely behavior in the form of the vector p of priors. We now ask whether and when it would in fact be rational for the players, having formulated p, to rely on p for all the information needed in determining what it is rational to actually do during the game.

As indicated above, the expectations of player i concerning the other players' strategies are given by the i-incomplete prior $p_{-i}$. Thus it might seem that each player i will play a strategy $q_i^o$ that is his <u>best reply</u> to this prior $p_{-i}$. Formally, he will determine a strategy $q_i = q_i^o$ that solves

$$\text{MAX} \quad H_i = H(q_i, p_{-i}) \tag{2.1}$$

where $H(\cdot)$ is i's payoff (utility) function. Generalizing this train of thought, it might seem that the solution to the game would be the vector $q^o = (q_1^o, \ldots, q_n^o)$ consisting of the set of these best replies for all n players. HS call this simplistic, first-order solution the <u>naive Bayesian approach</u>.

Why is this straightforward solution theory "naive"? The answer is simple: $q^o$ will not in general be an equilibrium point of the game G. Player i's strategy $q_i^o$ will not be a best reply to player j's strategy $q_j^o$, etc. As is well-known, in order for a strategy n-tuple to qualify as a bona fide solution in game theory, it must constitute an equilibrium point of a game.

It should prove fruitful to examine in more detail the ostensible

conflict between "naive Bayesian rationality" on the one hand and

game theoretic (equilibrium point) rationality on the other hand.

In analyzing this conflict, we shall see not only the key idea

underlying the tracing procedure, but also exactly why it is that

HS speak of their new theory as Bayesian in spirit.

One of the central tenets in Bayesian statistics and decision theory

is that people start off at some point in time  t  with information

called their prior (or "first-order") information. This initial information

can be denoted $I^P_t$.  Then as we move to time periods  t+1, t+2...,t+k,

they acquire new information and hence "update" their prior information

$I^P_t$ with the new information ("second-order") they obtain at each

subsequent point in time. They end up with $I^P_{t+k}$. In decision theory

proper, the proper way to "update" one's first order information with

second-order information is via Bayes rule of probability theory.

An analogue of this information updating procedure lies at the heart

of HS's TP. It is for this reason that their approach is called

Bayesian (also see 2.3 below). In the TP, the vector  p  of priors

is regarded to be the players' first-order information. Any subsequent

information they might receive is called second-order information.

But why would there be any need for second-order information? After

all, the given game G, since it is assumed to be well-defined,

should contain all the information that is necessary and relevant, or

should it? Let us press this question further.

HS distinguish between the specific information that is contained in the vector $p$ — first order information, and any additional information that the players may come to have about the likely reactions of one another to the strategy implications implicit in their first order information $p$. This latter type of information is the second-order information mentioned above.

To see all this more clearly, let us look at two distinct cases. First, suppose that the strategy n-tuple $q^o$ obtained when each player solves (2.1) above does, upon reflection, turn out to be "rational" for all n players. Then this conclusion about strategic rationality is regarded as second-order information in the HS theory. This would be the case were $q^o$ in fact an equilibrium point of the given game G.

Next suppose that $q^o$ is not in fact an equilibrium point, and this will generally be the case. If $q^o$ is not an equilibrium point, then by definition there must be at least one player i whose naive Bayesian strategy $q_o^i$ is not a best reply to the strategy (n-1) tuple $q_{-i}^o$ of the other (n-1) players. When the players realize this, they will realize that the naive Bayesian strategy $q^o$ is "recommending" irrational behavior on the part of at least one player. More specifically, second-order information about strategic rationality has come into conflict with the first-order information summarized by $p$.

How is this conflict resolved? HS make the generalized Bayesian

assumption that in the event of such conflicts (e.g. when $q^o$ is

not an equilibrium point) the players take their naive Bayesian

information (and its implied strategic behavior) $p(q^o)$ as a

starting point. They then "update this information" with second-order,

third-order,...,$k^{th}$-order information as this information is obtained.

But exactly how does this updating take place? Is there a game

theoretic analogue to Bayes' theorem that is called upon at analogous

junctions in their theory? Yes: HS implicitly assert that the tracing

procedure (about to be formally defined) is the appropriate analogue

to Bayes rule. It serves to update $p$ by feeding in subsequent

information in a strategy updating process until such time as

there remains no discrepancy at all between first and second order

information.

In the following Section 2.2, the TP is formally introduced and

characterized. However, we wish to emphasize that the reader who

really wishes to understand how the TP works will wish to read Chapter

3 as well as Section 2.2 below.


## 2.2 The Linear Tracing Procedure

In Section 2.1 we have seen how the HS theory leads to the formation

of prior strategic information and naive Bayesian behavior.

At the beginning of their decision-making process, the players'

expectations (first-order information) about each other's behavior
are given by the prior  p while their tentative strategy plans
(first-order behavior) are given by the strategy n-tuple $q^o$. The
TP itself is best viewed as an operator that takes  p  as given
and gradually updates it by feeding back information about strategic
rationality. More specifically, it is a Bayesian operator precisely
because of its role that is analogous to the role of Bayes' operator.
As the TP operates, i.e. as updating takes place, both  p and $q^o$
are subjected to systematic and continuous transformations until
they  both finally converge on a specific  equilibrium point q* of
the given game G. Thus, at the end of the process, a solution q*
is obtained according to which there is complete agreement between
the players' actual strategy plans and their expectations.

We shall now introduce the most accessible and transparent version
of the TP. It is called the linear tracing procedure. It is always
used in the HS theory except when mathematical complications require
a more general procedure (the logarithmic TP) to be substituted.

The linear tracing procedure (henceforth simply TP) is based upon
a one-parameter family of games $\{G^t\}$ with $0 \leqslant t \leqslant 1$. This one-
parameter family of games is derived from the original given game
G as follows. In any game $G^t$, every player i (i = 1,...,n) will have
the same strategy set $Q_i = \{q_i\}$ as in the underlying game G. But
his payoff function $H_i^t$ in $G^t$ will be different, defined now as

$$H_i^t(q_i, q_{-i}) = t\, H_i\,(q_i, q_{-i}) + (1-t)\, H_i\,(q_i, p_{-i}) \qquad (2.2)$$

where $H_i$ is i's payoff function in the game G. When $t = 1$ we have

$$H_i^1\,(q_i, q_{-i}) = H_i\,(q_i, q_{-i})$$

so that $G^t = G$. Yet when $t = 0$ we have

$$H_i^0\,(q_i, q_{-i}) = H_i\,(q_i, p_{-i}).$$

Thus $G^0$ is a rather special game in which the payoff $H_i^0$ of each i will depend only upon his own strategy $q_i$ and will be independent of the other players' strategies $q_{-i}$. Thus $G^0$ decomposes into n mutually separate maximization problems, one for each player. It can be shown that in most all cases of interest herein, the separable game $G^0$ has one equilibrium point (which in fact is a strong equilibrium point in pure strategies); and that moreover this equilibrium point $q^0$ is identical to the best-reply combination $q^0$ prescribed by naive Bayesian rationality (recall (2.1) above).

In discussing the games $G^t$ and in comparing them with the original game G, it will be helpful to make use of the following terminology. Instead of saying that $q_i^*$ is a best reply to $q_{-i}$ in game $G^t$, we shall simply state that $q_i^*$ is a $G^t$-best reply to $q_{-i}$. Likewise, we shall simply say that a point q that is an equilibrium point of a game $G^t$ is a $G^t$-equilibrium point.

By John Nash's fundamental theorem of 1951, it is known that every well-defined non-cooperative game G has at least one equilibrium point. Thus for any game $G^t$ the set $E^t$ of all equilibrium points will be nonempty. Let $X = X(G,p)$ be the graph of the correspondence $t \to E^t$ where t lies between 0 and 1. X can be shown generally to consist of a collection of pieces of one-dimensional (algebraic) curves. Consider a point x of X. It will have the form $x = (t,q)$ where q is an equilibrium point of $G^t$. HS call t the t-coordinate of x, and they call q the strategy part of x. Clearly, since $t \epsilon I = (0,1)$ (the closed unit interval) whereas $q \epsilon Q$ is a copy of the strategy space of G (or equivalently of $G^t$), we see that the graph X will always be a subset of the closed cylindrical set $Y = I \times Q$.

Suppose that the graph X contains a path L connecting a point $x^0 = (0,q^0)$ corresponding to an equilibrium point $q^0$ of the separable game $G^0$, with a point $x^1 = (1, q*)$ corresponding to an equilibrium point of the original game $G^1 = G$. Then HS call L a feasible path while $x^0$ and $x^1$ will be called the starting point and the end point, respectively, of the path L. Moreover, the strategy part q* of this end point $x^1$ is called the solution of the original game G as given by the path L.

A formal definition of the TP is now possible. It simply selects and defines as the solution to G the point q* obtained by "tracing" (i.e. by following) a feasible path L from its starting point $x^0$ to its endpoint $x^1$. The fundamental theorem of HS in this regard is that under suitable conditions for the linear TP to apply, the TP will follow one path L*,

called the <u>distinguished path</u> of G. Since this path L depends on
both G and p we shall write L = L(G,p). Moreover, HS have proved
that for any given pair (G,p) for which the (linear) TP is well-
defined, the procedure will always select a <u>unique</u> equilibrium
point $q^* = (q_1^*,\ldots,q_n^*) = T(G,p)$. In short, the TP is usually well-defined
and it works.

This completes our abstract description of the TP. Before turning
to an application in Chapter 3, we shall conclude the present Chapter
with a decision theoretical and game theoretical interpretation of
the HS procedure. Hopefully, this will assist the reader in appraising
the conceptual soundness and the applicability of the TP.


## 2.3. <u>A Decision/Game Theoretical Interpretation of the TP</u>

In this concluding Section, we provide a brief decision theoretical
and game theoretical interpretation of the TP. In 2.3.A we raise
a couple of points comparing the TP with the Bayes operator in
decision theory. Then in 2.3.B we present HS's own interpretation.

### 2.3.A. <u>Comparison with the Bayes Operator</u>

Earlier in Section 2 we have mentioned the sense in which the
TP is analogous to the "updating" that goes on in statistics
and decision theory more generally via Bayes rule. A couple of
points might help flesh out this analogy. First, it should be noted
that Bayes Theorem is not a "model" that was conceived to model

real-world cognitive processes. Bayes Theorem simply follows from

the elementary axioms and rules of probability theory. Indeed, until

the primitives that are operated upon by Bayes Theorem

were given a certain epistemological interpretation by modern statisticians

and decision theorists, the formula had no particular meaning.

Once substance is ascribed to the primitives, however, the formula

then assumes a very real meaning: it shows how "prior" information can

be updated to include "new information". One can go even further.

If one accepts the Savage axioms of decision theory, then it can be

asserted that Bayes rule specifies how one should update prior information.

The TP is quite different. It does not follow from anything analogous to

the axioms of probability. It is intended to be, and admitted to be

a model of the cognitive convergence of players' conflicting expectations

and intentions. While we feel the TP is a highly original and significant

approach to the most fundamental problem in all of game theory, namely

that of selecting a "best" equilibrium point in a game, we also feel

HS should discuss where the TP might stand in a whole family of alternative

approaches to this same problem. Are there possible axioms sets that

would identify the TP as the only/best manner in which to bring about

convergence of strategic intentions and expectations? As of yet we

do not know. Related to this matter is the problem of the logical

status of the TP. HS point out it is simply a model. But is it ultimately

a normative model, a conditionally normative model, a descriptive model,

a positive model, or some combination of these?

A second salient difference between the Bayes operator and the TP
is to be found in the matter of the "finality" of the results that
each produces. In a statistical decision problem, there is in principle
no end to the amount of new information that is obtained (or could be
obtained) by the decision maker. Hence there is no "terminus" to
the operation of Bayes rule. It is applied forever, at least in principle.
In stating this, we are talking about the general class of (statistical)
decision problems, and not simply the "optimal stopping problem" which
is indeed somewhat different.  In contrast to this, there is a well-
defined endpoint to the operation of the TP. When the end of the distinguished
path $L^*$ is reached, and a strategy n-tuple $q^*$ is identified, the
process is finished. The  discrepancy  between the prior information and
expectations and the actual expectations and information at $t = 1$ has
disappeared, as explained in 2.2. This distinction may be of interest
if and when researchers investigate entire classes of updating processes.

### 2.3.B. HS's Decision/Game Theoretical Interpretation

HS view the TP as a model of a solution process. It is seen as a
process of analysis, computation, expectation formation, and strategy
choice. The process is one in which n players will come to adopt
-- and will come to expect each other to adopt -- one particular
strategy n-tuple $q^*$ as the solution to their decision problem.
As time passes, the process  models gradual changes that take place
in i's own tentative strategy plans and in his expectations about
the other players' likely strategies. At any given moment, his plans

and his expectations will be determined by a specific point $x^t = (t,q)$ of the path $L*$. HS call this point $x^t$ i's <u>position point</u> at time t. As the process operates, $x^t$ will continuously move along $L*$ starting at $x^0$ and ending at $x^1$.

At any given moment, where i's position point is $x^t$, his tentative strategy plan will be to use the strategy $q_i$ prescribed by the strategy part $q = (q_1,..,q_i,..,q_n)$ of this point $x^t$. But he will know that as long as t is less than unity, the strategy part q of $x^t$ may not in fact correctly indicate the final solution of the game. Indeed, with t less than unity, q will in general not even be an equilibrium point of the given game G. Consequently, i will know that, even if $q_j$ is the strategy preferred by player j at the moment, there is no assurance either on his part or anyone else's part that $q_j$ will in fact be j's <u>final</u> and actual strategy.

More formally, HS argue that i will assign only probability t to the hypothesis that the other (n-1) players' behavior will in fact correspond to the strategy combination $q_{-i}$ predicted by the strategy part q of the position point $x^t$. He will reserve the remaining probability (1-t) to his <u>original</u> hypothesis -- the naive Bayesian hypothesis -- that the players' behavior will correspond to the original (naive) behavior $p_{-i}$. As HS further show, this implies that his actual expectations about the other players' behavior will correspond to the probability distribution

$$r_i(t,q) = t(q_{-i}) + t(q_{-i}) + (1-t)(p_{-i}).$$

This in turn can be used to justify the assumption that i's tentative

strategy plan at that moment will be to use strategy $q_i$. This is

so because it can be shown that $q_i$ will be i's best reply to the

probability distribution $r_i$ which, as just seen, represents his

expectations about the others' behavior.

Thus at each moment of the solution process, the t-value associated

with the position point $x^t$ will measure the degree of confidence

that each player i has in the tentative prediction provided by this

point $x^t$, namely that the other players will use strategy $q_{-i}$ as

prescribed by the strategy part q of $x^t$. As the process moves from

$t = 0$ to $t = 1$ all players will gradually move from a state of

complete predictive uncertainty to a state of complete confidence

in the behavior that will in fact be adopted.

At the start of the process, the complete uncertainty of i about

$j \neq i$'s behavior expresses itself in i's exclusive reliance on

the prior $p_{-i}$. He uses this prior in forming his expectation in

that he will entertain a probability distribution $r_i(0, q^o)$ assigning

probability 1 to $p_{-i}$ while assigning probability 0 to the strategy

combination $q_{-i}o$ associated with his initial position point $x^0$.

In contrast to this, at the end of the process, he will be in a

state of full predictive certainty because he will feel able to predict

which strategies the others will use. This is seen in the fact that

his probability distribution $r_i(1, q*)$ assigns unit probability to

the $(n-1)$ tuple $q_{-i}^*$ prescribed by the solution $q^*$ to G. In terms

of HS's distinction between first and second-order information,

at any position point $x^t = (t,q)$ the first-order information to which

each player i will react consists of the probability distribution

$r_i(t,q) = t(q_{-i}) + (1-t)(p_{-i})$. In contrast to this, his second-order

information consists of knowing that the others' tentative reactions

to their own first-order information will be the strategy $(n-1)$-tuple

$q_{-i}$. Formally, this second-order information can be defined as the

probability distribution $(q_{-i})$ generated by $q_{-i}$.

## 3. A WORKED APPLICATION

In the previous chapter, an analytical overview of the HS TP was
given. As promised, we now turn to a practical application of
the TP. In Section 3.1 we describe the game to be solved. In 3.2
we discuss the question of how the prior probability vector p is
determined, and we also discuss the use of the TP in defining certain
"risk dominance" relationships that hold among different equilibrium
points. In 3.3 we turn to numerical computations, and characterize
the solution given by the TP. 3.4 is an interpretative conclusion.

### 3.1. A Two-Person Bargaining Game with Incomplete Information

Drawing on the recent results of HS, we shall analyze a one-parameter
family of two-person bargaining games with incomplete information,
namely G(w), where  w  is the parameter that is to be varied. Our
ultimate purpose is to see what "solution" the TP yields for this
family of games as  the parameter  w  is varied. Let us now substantively
describe this family G(w) of games.

In any given game G(w), there are two players I and II who must decide
upon how to divide $100. Both players are assumed to have linear
utility functions for money. For convenience, assume that these utility
functions are normalized so that they ascribe  x  units of utility
to $x. In the event that the players cannot reach agreement on how

to divide the money, they receive the conflict payoffs $c_I$ and $c_{II}$.

We assume that $c_I = 0$ whereas $c_{II} = 0$ <u>or</u> w with w lying in the closed

interval ($0, $50). Whether $c_{II}$ is in fact equal to 0 or to w is

determined by a chance move at the start of the game. More specifically

these two possible conflict payoffs for player II each have probability

1/2. This probability distribution is assumed known to both players.

However, the true value of $c_{II}$ is actually only known to player II.

At this point, two observations are in order. First, the class of

bargaining games G(w) we are analyzing is a subset of the class

of <u>games with incomplete information</u>. This is so because, by definition,

a game possesses incomplete information if at least one player does

not know the utility payoffs of at least one antagonist. In the

present case, player I does not know whether his opponent's (conflict)

payoff is 0 or w. Second, it is now clear what the "parameter" is

in the one-parameter family of games G(w) we are analyzing: namely,

II's conflict payoff. For simplicity, we shall describe II as

being of a <u>weak type</u> or of a <u>strong type</u> -- i.e. being in a weak

or in a strong position -- according as whether $c_{II} = 0$ or w.

By the end of the chapter we shall have seen how the <u>solution</u>

prescribed by the TP changes as the <u>parameter</u> w changes*.

---

* HS themselves have investigated (and will publish) a more general
class of games than that looked at herein. In their work, w can
range from 0 all the way to 100. Guth and Selten (10) have analyzed
the same class of problem from a mathematically different point of view.

A game $G(w)$ is played in the following manner. Each player will

choose a number x lying between 0 and 100, to be interpreted as the

payoff proposed for player I. The number chosen by player I

(by player II) will be denoted $x_I$ ($x_{II}$). If $x_I = x_{II} = x$ then player

I will receive the payoff $u_I = x$ whereas player II will receive

the payoff $u_{II} = 100 - x$. On the other hand, in the event of conflict

where $x_I \neq x_{II}$, the two players will only receive their conflict

payoffs $c_I$ and $c_{II}$.

At this point it is appropriate to transform the game $G(w)$ described

above into what HS call the <u>agent normal form</u> of the game. To

understand this idea, recall that the "normal form" of any game

is simply a description of the game in terms of (i) the players,

(ii) the players' pure strategies, and (iii) the (utility) payoffs

corresponding to any choice of a strategy n-tuple by the players.

The "agent normal form" of a game is a straightforward generalization

of the normal form. Specifically, each "player" is replaced by

his "agents". It will be recalled from elementary game theory that

a given player has as many agents as he has information sets.

In our game, since player I has only one information set, he will

have only one agent. But player II, who has two information sets,

is replaced by two agents whom we shall call players 2 and 3.

In the sequel, Roman numerals refer to the two-player representation

of $G(w)$ whereas Arabic numerals refer to the three-player representation.

It traditionally has been thought that a game could be fully described by its normal form. However, HS show in their work that it is necessary to use the more general agent normal form if certain critical information is to be retained in passing from the extensive form to a normal form. The game theoretical issues that arise concerning this point will not be discussed in the present paper.

In $G(w)$, when player 1 selects a specific number $x_1 = x_I$ we shall say that he uses the pure strategy $s_1 = x_1$. Thus 1's pure strategy set is simply the closed interval $S_1 = (0,100)$. Likewise, a pure strategy $s_2$ of player 2 will be defined as a specific number $s_2 = x_2 = x_{II}$ that II may choose in the event that he is weak; and a pure strategy $s_3$ of player 3 will be defined as a specific number $s_3 = x_3 = x_{II}$ that player II may choose in the event he is strong. Thus 2's pure strategy set will likewise be the interval $(0,100)$. However 3's strategy set $S_3$ will be the smaller interval $(0,100-w)$. This is the case because if player 3 used any pure strategy $s_3 = x > 100-w$ then he would irrationally be accepting a payoff $u_3 = 100 - x < w$ which would be _less_ than the payoff $u_3 = c_3 = w$ he could obtain by provoking conflict. In what follows, the payoff function $u_i = U_i(s_1, s_2, s_3) = U(s)$ will be written as such, where $s = (s_1, s_2, s_3)$. Also, we shall write $u = (u_1, u_2, u_3)$ and $u = U(s)$.

HS show that in the class of games $G(w)$ under consideration, the

equilibrium points will consist of all strategy combinations of

the form $s = (s_1, s_2, s_3) = (x,x,x)$ with $0 \lessgtr x \lessgtr 100-w$. Any such

equilibrium point (EP) amounts to an agreement between all three

players, and will give a payoff vector $U(s) = (u_1, u_2, u_3) =$

$(x\ 100-x, 100-x)*$. What will be shown in the sequel is how the TP

selects a <u>unique</u> EP from this infinity of EPs, and moreover, how

this unique EP (the "solution") changes as the parameter $w$ varies.

But before proceeding with numerical computations, we shall discuss

the problem of generating the prior probabilities $p$, and we

shall introduce the concept of "risk dominance" as it arises in

the HS theory.

## 3.2. <u>Prior Probabilities and Risk Dominance in the HS Theory</u>

In Chapter 2 we discussed the role of the prior probability distribution

vector $p$ in the HS theory. Let us now discuss how these prior

probabilities are determined. Recall that a player's (e.g. player i's)

prior probabilities express the likelihood he assigns to the

adoption of different strategies at the <u>beginning</u> of the TP.

In other words, the strategies implied by $p$ constitute what we have

---

* There is in fact one other EP in a game $G(w)$, and this is the
strategy combination $s = (100,0,0)$ reflecting complete disagreement
between player I and II. However it can be shown that according to
the TP, the "full agreement" EPs dominate this full disagreement outcome.
We shall not discuss it further.

called "naive Bayesian behavior". More formally, for any player

i  (i = 1,2,3) in G(w),  let  $p_i = (p_i^1, p_i^2)$ be the probabilities

that are assigned to the strategy doublet* $s_i'$ and $s_i''$. Clearly

$p_i^1 + p_i^2 = 1$.

In the HS theory, $p_i$ is computed as follows. The theory makes the

prior probability that player i will use strategy $s_i'$ or strategy

$s_i''$ <u>proportional to his net payoffs</u> associated with the corres-

ponding equilibrium points s ' and s ", respectively. Let us see

how this happens at a more analytical level in the new theory.

Player i is chosen at random from a population $B_i$ of potential players

i, in which each potential player is characterized by a personal

parameter $b_i$ that is assumed uniformly distributed over the closed unit

interval (0,1).  Any player i, characterized by a specific value of

$b_i$, will assume that the other two players will either both use

their s'-strategies $s_j'$ and $s_k'$, or will both use their s"-strategies

$s_j''$ and $s_k''$ -- with the former event having probability $b_i$ whereas

the latter event has the probability $(1-b_i)$. Denote

$$U_i(s_i',s_j',s_k') = u_i' \text{ and } U_i(s_i'',s_j'',s_k'') = u_i'' \quad (3.1)$$

Also, by our assumptions about G(w) from Section 3.1, we have

---

* The reason why we are restricting ourselves here to a pair of
strategies (rather than to an n-tuple) will be made clear in
the sequel.

$$U_i(s_i'',s_j',s_k') = U_i(s_i',s_j'',s_k'') = c_i \qquad (3.2)$$

Thus, player i's <u>expected</u> payoff, if he uses strategy $s_i'$, will be

$$v_i' = b_i u_i' + (1 - b_i)c_i \qquad (3.3)$$

while if he uses strategy $s_i''$ then his expected payoff will be

$$v_i'' = b_i c_i + (1 - b_i) u_i'' \qquad (3.4)$$

By the expected utility principle of decision theory, player i will use strategy $s_i'$ if $v_i' > v_i''$, that is (by simple algebraic substitution and rearrangement of terms) if

$$b_i \geqslant \frac{u_i'' - c_i}{(u_i'-c_i) + (u_i''-c_i)} \qquad (3.5)$$

And if this inequality is reversed, he will use strategy $s_i''$. Since the HS theory assumes that the population $B_i$ is uniformly distributed over the unit interval as stated above, the probability that the former alternative will hold is clearly

$$p_i^1 = \frac{u_i' - c_i}{(u_i' - c_i) + (u_i'' - c_i)} \qquad (3.6)$$

whereas the probability that the latter alternative will hold is

$$p_i^2 = \frac{u_i'' - c_i}{(u_i' - c_i) + (u_i'' - c_i)} \qquad (3.7)$$

Here we see exactly how the HS theory makes the prior probabilities proportional to the net payoffs associated with the strategies.

One point that may trouble the reader at this junction is why HS
make the assumption about the uniformity of the probability dist-
ribution over $B_i$, the set of possible "types" of player i. This is
a matter that HS will explain in some detail in their forthcoming
book. For the moment, it will hopefully suffice to have shown
that _use_ of this assumption (in the above derivations) yields the
desired result: namely prior probabilities that are proportional to
the net payoffs associated with the different strategy alternatives.

As we know from Chapter 2, the purpose of the TP itself is to
determine a unique equilibrium point of the game as the "solution".
In the present context, what this amounts to is determining whether
strategy s' or s" is to be the solution. At this point we introduce
the notion of _risk dominance_. In comparing and contrasting s' and s"
via the theory of the TP, HS speak of the risk dominance relationships
between s' and s" that are revealed by the workings of the TP.
As we shall hopefully show, explicating the TP in terms of risk
dominance adds intuitive insight into what is really going on. Let
us now develop these points analytically in the remainder of 3.2.

Starting with the triplet $p = (p_1, p_2, p_3)$ of prior probability distributions,
the TP itself is applied to an auxiliary family of games $G(w)^t$
for $0 \leqslant t \leqslant 1$. (Henceforth, we shall simply write $G^t$ for $G(w)^t$.)
The payoff functions for the auxiliary games are given by

$$U_i^t(s_i, s_j, s_k) = (1-t)\, U_i(s_i, p_j, p_k) + t\, U_i(s_i, s_j, s_k) \qquad (3.8)$$

as we know from equation (2.2). Let the set of equilibrium points

of any game $G^t$ be denoted $E^t$. Let P be the graph of the mapping

$t \rightarrow E^t$. In each game $G^t$ , we know from Chapter 2 that the players

will use a strategy n-tuple $s^t$ corresponding to a distinguished path

L of the graph P. In the initial game $G^0$ corresponding to t = 0, this

distinguished equilibrium point strategy n-tuple $s^0 = (s_1^0, s_2^0, s_3^0)$

is computed as follows. The initial strategy $s_i = s_i^0$ of player i is

that strategy that maximizes his inital payoff, i.e.

$$U_i(s_i^0, p_j, p_k) = MAX\ (u_i'^0, u_i''^0) \qquad (3.9)$$

where

$$u_i'^0 = U_i(s_i', p_j, p_k) \qquad \text{and} \qquad u_i''^0 = U_i\ (s_i'', p_j, p_k) \qquad (3.10)$$

The maximization of $U_i$ can be restricted to the two strategies $s_i'$ and

$s_i''$ since all other strategies $s_i$ will drive $U_i$ to zero.

Let us now see what happens analytically as time goes on, i.e. as

t becomes greater than zero and the tracing procedure unfolds. We shall

be particularly interested in the conditions under which each

of the players _deviates_ from his initial  (naive Bayesian) strategy

choice: for it is in the analysis of these deviations that the notion

of risk dominance is made clear.

At t = 0, suppose that the players' initial strategies are $s_i^0, s_j^0, s_k^0$.

And more specifically, suppose that $s_i^0 = s_i'$. If the other two

players now go on using the same strategies $s_j^0$ and $s_k^0$ at higher t

values, then player i will not deviate from his initial strategy
$s_i^0 = s_i'$ as long as the following inequality is maintained:

$$U_i^t(s_i',s_j^0,s_k^0) > U_i^t (s_i'',s_j^0,s_k^0). \qquad (3.11)$$

Likewise if the initial strategy $s_i^0 = s_i''$, then player i will stick
to this as long as the converse inequality is preserved. In short,
player i will only deviate from his initial strategy at that point $t = t_i$
when

$$U_i^t(s_i',s_j^0,s_k^0) = U_i^t(s_i'',s_j^0,s_k^0) \qquad (3.12)$$

or equivalently at that point $t = t_i$ when

$$(1 - t_i) \, U_i(s_i',p_j,p_k) + t_i U_i(s_i',s_j^0,s_k^0)$$

$$\qquad (3.13)$$

$$= (1 - t_i) \, U_i (s_i'',p_j,p_k) + t_i U_i (s_i'',s_j^0,s_k^0)$$

This important equation (3.13) yields a linear equation for computing
the time (the $t_i$ value) of deviation by i, if such a time exists at all.
Applying like reasoning to players j and k, we can arrive at expressions
for $t_j$ and $t_k$. The conclusion we reach is that i will be the first
to deviate if

$$t_i < \text{MIN} (t_j,t_k) \qquad (3.14)$$

provided only that $t_i$ exists. We now apply the above machinery to
the specific one-parameter family of games G(w) and see how the TP
works in practice. The exercise will provide even greater insight
into the notions of risk dominance and deviation-from-strategy.

## 3.3 Numerical Derivation of the Solution for the Family of Games G(w)

We now couple the foregoing analytical results with the numerical

parameters of the family of games G(w) and show what the TP implies

as the solution (distribution of \$100 among the players). Let $s' = x'$

and $s'' = x''$ be two equilibrium points of G(w). Assume for convenience

that $x' > x''$. Define z and v as

$$z = (x'+x'')/2 \qquad \text{and} \qquad v = x' - z \qquad (3.15)$$

Transposing we can write

$$x' = z + v, \quad x'' = z - v, \quad \text{with} \quad v > 0 \qquad (3.16)$$

Given these algebraic expressions, given the parameters of the

family of games G(w), and given the analytical results of section 3.2,

we know that the prior probabilities associated with s' and s'' are

$$p_1^{\ 1} = x'/(x'+x'') = (z+v)/2z \qquad p_1^{\ 2} = (z-v)/2z \qquad (3.17)$$

$$p_2^{\ 1} = (100-z-v)/(200-2z) \qquad p_2^{\ 2} = (100-z+v)/(200-2z) \qquad (3.18)$$

$$p_3^{\ 1} = (100-w-z-v)/(200-2w-2z) \qquad p_3^{\ 2} = (100-w-z+v)/(200-2w-2z) \qquad (3.19)$$

Let us now determine the conditions under which the three players will

chose either strategy $s_i'$ or $s_i''$ at the <u>start</u> of the game. Thereafter

we analyze the deviation behavior of the players under various sets

of circumstances.

From (3.17) and from the analysis of Section 3.2 above, we see that

$$u_1^{,0} = U_1(s_1{}',p_2,p_3) = (z+v)\left( \frac{1}{2}\ \frac{100\text{-}z\text{-}v}{200\text{-}2z} + \frac{1}{2}\ \frac{100\text{-}w\text{-}z\text{-}v}{200\text{-}2w\text{-}2z} \right) \qquad (3.20)$$

whereas

$$u_1^{''0} = U_1(s_1'',p_2,p_3) = (z\text{-}v)\left( \frac{1}{2}\ \frac{100\text{-}z\text{+}v}{200\text{-}2z} + \frac{1}{2}\ \frac{100\text{-}w\text{-}z\text{+}v}{200\text{-}2w\text{-}2z} \right) \qquad (3.21)$$

Player 1 will select strategy $s_1^0 = s_1{}'$ as his initial strategy if $u_1^{,0} > u_1^{''0}$. But this is equivalent to the requirement that

$$4z^2 - (600\text{-}3w)z + (20{,}000 - 200w) > 0 \qquad (3.22)$$

This requirement will in turn be satisfied if

$$z < f \qquad (3.23)$$

where

$$f = (600\text{-}3w - \sqrt{40{,}000\text{-}400w\text{+}9w^2}) \ / \ 8 \qquad (3.24)$$

By like reasoning, player 1 will select $s_1^0 = s_1''$ when

$$z > f \qquad (3.25)$$

In the case of player 2, the above reasoning quickly yields that player 2 will choose $s_2^0 = s_2{}'$ as his initial strategy when $u_2^{,0} > u_2^{''0}$ which in turn will obtain when

$$z < 50.$$

He will choose $s_2^0 = s_2''$ if $z > 50$. Finally, player 3 will choose $s_3^0 = s_3'$ when

$$z < 50 - (w/2) \qquad (3.26)$$

whereas he will select $s_3^0 = s_3''$ when (3.26) is reversed.

Now, consider two important points about the quantity $f$ originally defined in (3.24). First, it is seen that

$$50 - (w/2) < f < 50 \qquad (3.27)$$

Second, solving (3.24) gives the following values of $f$ as a function of the conflict payoff parameter $w$ (at \$10 intervals):

$$
\begin{array}{lcccccc}
w: & 0 & 10 & 20 & 30 & 40 & 50 \\
f: & 50 & 47 & 43.9 & 40 & 35.5 & 30.1
\end{array}
\qquad (3.28)
$$

Drawing upon the earlier results of Section 3.2 and upon equations (3.22)-(3,28) it will now be seen that <u>the equilibrium point</u> $s^* = f(w)$ <u>is the solution to our one-parameter family of games</u>. Or, to use the language of the TP, it will be seen that this particular equilibrium point risk dominates all other possible solutions.

The above results yield four possible sets of possible interrelationships between the parameters $w$, $f$, and $z$. Let us now see what the TP implies in each of these possible cases. Specifically, let us analyze which players deviate from their <u>initial</u> strategies (determined just above) and what this implies in each of the four cases.

Case 1: $z < 50 - (w/2)$: in this situation, $s^0 = (s_1{}^0, s_2{}^0, s_3{}^0) =$
$(s_1', s_2', s_3') = s'$. $s'$ is, we recall, a (strong) equilibrium point
in the given game $G(w)$. But when the players' initial strategies
constitute a strong equilibrium point of the given game, then
according to a HS lemma, they will go on using that strategy n-tuple
during the duration of the TP. Hence in Case 1, $s'$ risk dominates $s''$.

Case 2: $50 - (w/2) < z < f$. In this event, $s^0 = (s_1{}^0, s_2{}^0, s_3{}^0) =$
$(s_1', s_2', s_3'')$. This strategy combination is clearly not an equilibrium
point in the given game $G(w)$; thus at least one player will deviate
from his original strategy during the operation of the TP. Applying
the theory and the formulas developed in Section 3.2, we find out
that it is player 3 who will first deviate from his initial
strategy $s_3{}^0 = s_3''$, thereupon adopting $s_3'$. The strategy triplet
in effect at that switching time is -- in the same manner as we
saw just above -- an equilibrium point of the given game. Hence it
is followed until $t = 1$. So once again $s'$ risk dominates $s''$.

Case 3: $f < z < 50$. Here $s^0 = (s_1'', s_2', s_3'')$. This is clearly not
an equilibrium point of the given game. Via some rather complex
calculations, HS show that player 3 has no incentive to deviate
from his initial strategy. What is at issue is whether player 1
or player 2 will deviate first, i.e. the values of $t_1$ and $t_2$
as defined in Section 3.2. It turns out that who deviates first
depends upon the numerical values of $x'$ and $x''$ and also upon the
value of $w$. However, what holds generally is that if $t_1 > t_2$
then $s'$ risk dominates $s''$ whereas the converse holds if $t_1 < t_2$.

Case 4: $z > 50$. In this case, $s^0 = (s_1'', s_2'', s_3'') = s''$. This is of

course an equilibrium point of the given game, so we see that $s''$

risk dominates $s'$. The set of inequalities we have observed yields

the desired result (among others), namely that the equilibrium

point $s^* = (f,f,f)$ risk dominates all other equilibrium points.


## 3.4 An Intuitive Interpretation of the Solution to G(w)

At this point let us look at the numerical solution to $G(w)$ given

above by the TP, and let us examine how "reasonable" it is.

In investigating this matter, we shall refer back to the classical

Nash solution to the two-person bargaining problem referred to in

Chapter 1. And in so doing, we shall compare the TP solution

to $G(w)$ with the Nash solution to a game that approximates

$G(w)$.

Let us first reproduce the numerical solution as originally given in
(3.28):

| w: | 0 | 10 | 20 | 30 | 40 | 50 |
|---|---|---|---|---|---|---|
| f: | 50 | 47.2 | 43.9 | 40 | 35.5 | 30.1 |

(3.28)

The solution $s^* = (f,f,f)$ tells us that all three players will

agree that player 1 should receive $f(w)$ whereas player two (or

three, depending on II's true type) will receive $100 - f(w)$.

Thus, we see from (3.28) that if $w = 0$, each player gets $50,

if $w = 10$, player 1 gets $47.20 whereas player II gets $52.80, and

so on. In sum, the higher the conflict payoff to player II, i.e.

the higher the value of w, then the lower the final payoff of money
to player 1 whose conflict payoff is known to be 0. The TP yields
an intuitively appealing result that corresponds to our sense about
the ability of a "stronger" player to bargain down a "weaker" player
-- where strength/weakness is given by the conflict payoff of
player II.

Let us see how the TP solution relates to the classical Nash
solution to the extent that such a comparison is possible. (Recall
that Nash did not define a solution for two player games with
incomplete information, so that strictly speaking he did not develop
a theory that could solve $G(w)$.) Understanding the relation
between the TP solution and the classical Nash solution is par-
ticularly insightful since it was in analyzing the Nash solution
that John Harsanyi (following F. Zeuthen) first developed his
original theory of risk dominance.

In $G(w)$ player 1 has an equal chance of facing player 2 or player
3 as his antagonist. Thus, as HS point out, our game $G(w)$ lies
half way between a two-person game $G'$ in which player 1 would
with certainty face player 2 as an opponent, and a game $G''$
in which he would with certainty face player 3. By applying the
classical Nash solution (or equivalently the new TP) to the
true two-person game $G'$ we find that each player receives $50.
The intuitive reason for this of course is that $G'$ would be

a fully symmetrical game: both players have \$0 as their conflict
payoffs, and their utility functions are identical since each
player is assumed to have a function linear in money. Thus the
Nash solution calls for the players to split the gains equally.
Formally, $u_1' = u_2' = 50$. In $G''$ players 1 and 3 receive
$u_1'' = 50 - (w/2)$ and $u_3'' = 50 + (w/2)$. If we then think of the
original game $G$ as being a kind of "half-way house" between $G'$
and $G''$, it would seem that the players might receive as their
payoffs

$$u_1^* = (u_1' + u_1'')/2 = 50 - (w/4)$$

$$u_2^* = u_3^* = 100 - u_1^* = 50 + (w/4).$$

Of course, the original game $G$ is not in any formal sense a
half-way house between $G'$ and $G''$. Rather it is a true two-player
game with incomplete information. And as stated just above, the
classical Nash theory did not deal with these games. Nonetheless
it is interesting that in the one-parameter family $G(w)$ of games
we are considering, the TP solution is very close to the
hypothetical payoffs $u_1^*$, $u_2^*$, and $u_3^*$ exhibited above. Specifically,
we can see that

$$u_1(s^*) = f < u_1^* \quad \text{and} \quad u_2(s^*) = u_3(s^*) = 100 - f > u_2^* = u_3^*.$$

By simply substituting in various values of $f(w)$ given in (3.28)
the reader will immediately see that the TP results are in fact
very close to the generalized Nashian results.

We feel that this result is important because of the well-known
intuitive appeal of the Nash bargaining solution. Furthermore, it
should be pointed out that if we apply the tracing procedure to
any two-person bargaining game of <u>complete</u> information, we obtain
as the solution to the game the classical Nash solution itself.
What is particularly interesting to us about this fact is that
while the classical theory and the new theory make use of quite
<u>different</u> cognitive and psychological models (recall Section 1.2
of Chapter 1 above), both give <u>identical</u> results in the class of
games for which the classical theory is well-defined.

### 3.5: <u>Concluding Comments</u>

In this Chapter we have analyzed a particular class of two-person
games, namely two-person fixed threat bargaining games of incomplete
information. Choice of this particular example has proven useful for
several reasons. First, it has been possible to show how the tracing
procedure works with a minimum of technical detail. In particular,
it was not necessary to enter into a discussion of the concepts of
the <u>cells</u> and <u>formations</u> of a game. Had we chosen a more complex
example, this would have been necessary, and would have greatly
lengthened and complicated the present paper. Second, the example
permitted us to see how Harsanyi-Selten can take a game of incomplete
information and represent it as a well-defined game of complete
information by using their agent normal game form. Third, we have

just seen how our example permitted a link to be established between

the new tracing procedure theory and classical game theory. Thus

we feel we have achieved our mission as stated in Chapter 1: to

present a streamlined version of the tracing procedure theory that

emphasizes its basic ideas, and to show how these ideas work in practice

in the context of a concrete example. As pointed out in Section 1.4,

a price paid for this particular mission is the omission of many

technical considerations of great conceptual interest, in particular

the notion of cells and formations, and the concept of a uniformly

perturbed agent normal form. Nonetheless, the reader who has completed

this research report should have gained a good understanding of

the basic ideas of the new theory, their applicability, and their

relationship to the only earlier determinate solution theory, namely

the bargaining equilibrium model discussed in Section 1.3 of Chapter 1.

If he has gained this, we have succeeded in our efforts.

## BIBLIOGRAPHICAL REFERENCES CITED

1. Brock, Horace W., "A Critical Discussion of the Work of John C. Harsanyi", _Theory and Decision_, Volume 9, 1978.

2. Harsanyi, John C., "Approaches to the Bargaining Problem Before and After the Theory of Games", _Econometrica_, Volume 24, 1956.

3. Harsanyi, John C., "The Tracing Procedure: A Bayesian Approach to defining a Solution for n-Person Noncooperative Games", _International Journal of Game Theory_, Volume 4, 1975.

4. Harsanyi, John C., "A Solution Concept for n-Person Noncooperative Games", _International Journal of Game Theory_, Volume 5, 1977.

5. Harsanyi, John C., _Rational Behavior and Bargaining Equilibrium in Games and Social Situations_, Cambridge University Press, Cambridge, 1977.

6. Harsanyi, John C., "Analysis of a Family of Two-Person Bargaining Games with Incomplete Information", Working Paper #CP-604, Center for Research in Management Science, University of California at Berkeley (Unpublished), 1979.

7. Luce, R. Duncan, and Raiffa, Howard, _Games and Decisions_, John Wiley and Sons, New York, 1957.

8. Nash, John F., Jr., "The Bargaining Problem", _Econometrica_, Volume 18, 1950.

9. Nash, John F., Jr., "Non-Cooperative Games", _Annals of Mathematics_ Volume 54, 1951.

10. Selten, Reinhard, and Guth, Werner, "Macht Einigkeit Stark? - Spieltheoretische Analyse einer Verhandlungssituation", in _Neuere Entwicklungen in der Wirtschaftswissenschaften_, Duncker and Humblot, Berlin, 1978.

11. Zeuthen, Frederik, _Problems of Monopoly and Economic Warfare_, G. Routledge and Sons, London, 1930.

## BIBLIOGRAPHICAL NOTE

There were three principal sources underlying the present research effort: (i) many helpful discussions with John Harsanyi during 1978-1980; (ii) unpublished chapters from the forthcoming Harsanyi-Selten book A Noncooperative Solution Concept with Cooperative Applications (kindly given to the author by John Harsanyi in 1980); and (iii) reference (6) above -- the bargaining game application that will appear in its final form in the forthcoming Harsanyi-Selten book.

LMED 8